

# The Importance of Lexicalized Syntax Models for Natural Language Generation Tasks

Hal Daumé III, Kevin Knight, Irene Langkilde-Geary, Daniel Marcu and Kenji Yamada

Information Sciences Institute  
Computer Science Department  
University of Southern California  
4676 Admiralty Way, Suite 1001  
Marina del Rey, CA 90292

{hdaume,knight,ilangkil,marcu,kyamada}@isi.edu

## Abstract

The parsing community has long recognized the importance of lexicalized models of syntax. By contrast, these models do not appear to have had an impact on the statistical NLG community. To prove their importance in NLG, we show that a lexicalized model of syntax improves the performance of a statistical text compression system, and show results that suggest it would also improve the performances of an MT application and a pure natural language generation system.

## 1 Introduction

We distinguish between three types of language models:

- $n$ -gram language models look only at word sequences to gauge the quality of a sentence.
- Non-lexicalized syntax models consider only syntactic structure down to the level of word tags in assessing the grammaticality of a sentence. A PCFG is an example of such a model.
- Lexicalized syntax models take into account both sentence syntax and lexical values when determining the quality of a sentence.

The parsing community has long recognized the importance of lexicalized models of syntax for building robust natural language applications. For example, Charniak (1997) showed that by using a lexicalized syntax model instead of a non-lexicalized PCFG syntax model trained on Penn

Treebank data, one can increase the performance of a syntactic parser from 73.75% labeled recall and precision to 83.75%. More sophisticated lexicalized models of syntax (Collins, 1997; Charniak, 2000) have increased the performance of syntactic parsers to 90% labeled recall and precision. Lexicalized models of syntax have been also proven useful in speech recognition (Chelba and Jelinek, 1998) and language modeling (Charniak, 2001; Roark, 2001).

By contrast, lexicalized models of syntax do not appear to have had an impact on the statistical NLG community. Langkilde and Knight (1998), for example, use an  $n$ -gram model to select between different lexical renderings of a meaning representation. Knight and Marcu (2000) use a combination of bigram and context free probabilities to select between sentence compressions. To our knowledge, the only NLG work that resonates with the work in parsing, speech recognition, and language modeling is that of Bangalore and Rambow (2000), who show that a statistical generation system that uses a lexicalized hierarchical model of syntax outperforms a system that uses a random model.

Given the small interest in exploiting lexicalized models of syntax in NLG, we may conclude that such models have no role to play in this area. In this paper, we show that this is not the case. To prove the importance of lexicalized models of syntax in NLG, we focus on three distinct tasks, each involving a generation component. First, we show that a lexicalized model of syntax can improve the performance of a statistics-based text compression system. Second, we show that a lexicalized model of syntax may improve the outputs of a Chinese-to-English machine translation system. Finally, we analyze the results of a pure natural language genera-

Report Documentation Page			Form Approved OMB No. 0704-0188		
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>2002</b>		2. REPORT TYPE		3. DATES COVERED <b>00-00-2002 to 00-00-2002</b>	
4. TITLE AND SUBTITLE <b>The Importance of Lexicalized Syntax Models for Natural Language Generation Tasks</b>			5a. CONTRACT NUMBER		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S)			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>University of California, Information Sciences Institute ,4676 Admiralty Way, Marina del Rey, CA, 90292</b>			8. PERFORMING ORGANIZATION REPORT NUMBER		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)			10. SPONSOR/MONITOR'S ACRONYM(S)		
			11. SPONSOR/MONITOR'S REPORT NUMBER(S)		
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES <b>8</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

tion system.<sup>1</sup>

In each experiment, we assess the impact that a lexicalized model of syntax may have on improving the quality of existing systems using an off-the-shelf component: the parser built by Charniak (2000).

Our use of Charniak’s parser provides for a loose coupling of a lexicalized model of syntax with a generation system, which is far from ideal. Nevertheless, these experiments provide evidence that lexicalized models of syntax can improve the quality of the outputs of statistics-based generators. Our results motivate work aimed at building generation systems that choose between possible renderings of the same meaning using not only *n*-gram probabilistic models, but lexicalized models of syntax, such as those proposed by Collins (1997) and Charniak (2000).

## 2 Statistical Summarization

### 2.1 Document Compression

To assess the impact that lexicalized models of syntax may have on the task of summary generation, we used a noisy-channel document compression system (Daumé III and Marcu, 2002), which generalizes the sentence compression system developed by Knight and Marcu (2000) to the document level.

In Daumé and Marcu’s system, possible document compressions are packed into a shared-forest structure which contains explicit channel-model probability scores. These channel probabilities are based on discourse PCFG probabilities and syntactic PCFG probabilities, as well as compression probabilities. None of these probabilities are lexicalized.

We then use a generic forest ranker which combines these channel probabilities with bigram-based source model probabilities (Langkilde, 2000) to extract the top scoring compression. The only point at which lexicalization enters the model is in the bigram-based source model. This leads to the generation of many poor syntactic structures. For instance, the model does not know the difference between transitive and intransitive verbs, and therefore will often “compress off” the objects of transitive verbs.

---

<sup>1</sup>We are grateful to Eugene Charniak for suggesting these experiments to us.

After the system was constructed, we evaluated its performance on 16 documents chosen from the Wall Street Journal portion of the Penn Treebank, each containing between 41 and 87 words<sup>2</sup>. We also evaluated it on 5 documents selected from the Mitre corpus (Hirschman et al., 1999), each document containing between 64 and 91 words. We used two corpora to see whether the system’s performance varied with text genre. In the evaluations, the system outperformed a baseline system presented by Knight & Marcu (2000), applied iteratively. However, it still made many errors which would not have been made, had it had a good grasp of English grammar. We performed an exhaustive error analysis on the system to see where it could be improved.

### 2.2 Error Analysis

After analyzing the errors the system made, we classified them into ten classes, each listed in Table 1 with an example of the error and an example of this error fixed.

We then tabulated the frequency with which these error occurred in the summaries produced by the system. They are shown in Table 2.

Three of these error classes can be considered grammaticality errors: DET, MOD, COMP, NOVERB and NUM. These three alone account for 37 errors out of a total of 57 errors. It seemed that the integration of a lexicalized model of syntax into the source model would easily remove all of these problems. Since the document compression system was able to output an *n*-best list of possible compressions, we were able to rerank this list using Charniak’s parser.

### 2.3 Syntax-Based Reranking

To do this reranking, we modified Charniak’s (2000) parser so that instead of outputting the optimal parse tree for a given sentence, it outputs its non-normalized maximum entropy score. We then ran the modified parser on the 1000 best compressions according to the bigram model, normalized the probabilities by length and chose the single best compression.

---

<sup>2</sup>Because there are an exponential number of summaries that can be generated for any text, the decoder runs out of memory for longer documents; therefore, we selected shorter subtexts from the original documents.

## Grammatical Errors

**DET** - A noun or noun phrase is missing a determiner.

Erroneous “El Paso owns and operates refinery.”  
Fixed “El Paso owns and operates *a* refinery.”

**COMP** - The complement of a verb or noun is missing, rendering the output ungrammatical.

Erroneous “Banco Exterior was run by politicians who lacked the skills or the will.”  
Fixed “Banco Exterior was run by politicians who lacked the skills or the will *to do* . . . .”

**NUM** - Often if the original document contains a percentage (for instance “5%”), the number will be dropped without the percentage sign.

Erroneous “The rate is %.”  
Fixed “The rate is 5%.”

**NOVERB** - Sentence lacking a verb.

Erroneous “Stewart, the builder.”  
Fixed “*It is named for* Stewart, the builder.”

## Discourse/Coherence-specific Errors

**ANTE** - The antecedent of an anaphor has been dropped, resulting in incoherence.

Erroneous “Terms are subject to change, the company said.”  
Fixed “Terms are subject to change, Banco Exterior said.”

**CUE+** - Uninterpretable cue-word. For instance, if  $D_1$  and  $D_2$  form a discourse constituent and  $D_2$  is contrasting  $D_1$  and  $D_2$  begins with “But”, but  $D_1$  is dropped, so should be “But.”

Erroneous “But the proposed transaction calls for an exchange of the debt . . . .”  
Fixed “The proposed transaction calls for an exchange of the debt . . . .”

**CUE-** - A cue word or phrase is missing, rendering the output incoherent.

Erroneous “The President of the United States urged the armed forces to advance. His commanders did not have the initiative.”  
Fixed “President of the United States urged the armed forces to advance. *When he did*, his commanders did not have the initiative.”

## Summarization-specific Errors

**MOD** - A nominal modifier which should not have been dropped has been and significant meaning is lost.

Erroneous “Tons will fill damp barns across the land.”  
Fixed “Tons *of vegetables* . . . will fill damp barns across the land.”

**MISS** - The compression misses important information.

**EXTRA** - The compression contains unimportant information.

Table 1: Types of errors in the texts generated by the model.

This is far from an ideal language model. For instance, there is no guarantee that the structure the parser will derive for the sentence will be the same structure the compression model generated. Furthermore, since the parser assumes only one input sentence, the scores produced by the best parse for two different sentences may not be comparable.

## 2.4 Evaluation

After reranking, we performed the same error analysis as before. The results for the new error analysis are summarized in Table 3.

We can see from the delta row in Table 3 (negative

numbers are good), that we removed most grammaticality errors. We saw modest improvement with respect to the dropping of important modifiers; this is to be expected, though, since the syntax model has no idea of “importance.” The same argument explains the minimal change in the missing antecedent problems. We removed all instances of extra information but added seven additional counts of missing information.

Any summarization system must balance length of summary against informational and grammatical quality. In order to have more documents of higher grammaticality, more words are often nec-

	Grammaticality					Discourse			Summarization	
	DET	MOD	COMP	NUM	NOVERB	ANTE	CUE+	CUE-	MISS	EXTRA
wsj_0607	1	1								
wsj_0616	1		1						1	1
wsj_0632									1	
wsj_0654			1			1	1			
wsj_0655			1							
wsj_0667	1	2		1						
wsj_0689			1			1			1	1
wsj_1126		2					1		1	
wsj_1146		1								
wsj_1189	1	2			1					
wsj_1307	1	1	2					1		
wsj_1331									1	
wsj_1346			2						1	
wsj_1376		1	1					1		
wsj_1380	1		1							
wsj_2386			1						1	1
rm5-10	1					1			1	
rm5-22	1		1			2				
rm5-27		1	1							
rm5-6	1				1					
rm5-9										
<b>Count</b>	10	11	13	1	2	5	2	2	8	3

Table 2: Tabulation of the errors in the texts generated by the model.

	Grammaticality					Discourse			Summarization	
	DET	MOD	COMP	NUM	NOVERB	ANTE	CUE+	CUE-	MISS	EXTRA
wsj_0607										
wsj_0616		1							1	
wsj_0632			1						1	
wsj_0654									1	
wsj_0655		1						1	1	
wsj_0667				1						
wsj_0689		2						1	1	
wsj_1126									1	
wsj_1146		1								
wsj_1189										
wsj_1307									1	
wsj_1331		1							1	
wsj_1346		1								
wsj_1376								1	1	
wsj_1380		1						1		
wsj_2386									1	
rm5-10	1								1	
rm5-22			1						1	
rm5-27			1						2	
rm5-6										
rm5-9			1						1	
<b>Count</b>	1	8	4	1	0	4	0	0	15	0
<b>Delta</b>	-9	-3	-9	0	-2	-1	-2	-2	+7	-3

Table 3: Tabulation of the errors in the texts generated by the model with syntax-based rescoring.

essary, which reduces the amount of information which can be packed into a summary of comparable length. Here, by removing most of the grammaticality errors, we caused the system to drop some of the important information.

To determine whether the changes in the system

are noticeable to a user, we carried out a subjective evaluation. We presented 8 human judges with the outputs generated by the original text compression system, the results after rescoring, and human-generated compressions. These judges were asked to rank outputs on a scale from 1 to 5 (5 being the

	WSJ Texts				Mitre Texts			
	Cmp	Grm	Coh	Qual	Cmp	Grm	Coh	Qual
Old	0.47	3.11	2.98	2.55	0.47	3.57	2.90	2.80
Rescored	0.42	3.41	3.11	2.64	0.30	3.00	3.05	2.23
Hand	0.59	4.38	4.33	3.97	0.46	4.70	4.45	4.10

Table 4: Evaluation Results

best) on metrics of compression rate (Cmp), Grammaticality (Grm), Coherence (Coh) and Compression Quality (Qual). The results of this evaluation are summarized in Table 4.

In the Wall Street Journal data, there was a moderate improvement in grammaticality, coherence and quality, as error analysis suggested. In the Mitre data, grammaticality and quality went down significantly, while coherence remained steady. This can be attributed to two factors. First, there were few errors in the Mitre data to start with and thus less room for improvement; Second, the Mitre data is out of domain for both the document compression system and for the parser, which leads to less reliable statistics.

### 3 Machine Translation

For our machine translation experiments, we use the statistical MT system of Yamada and Knight (2001). This system produces English translations of foreign language sentences by exploiting three components:

- A Translation model (TM). For any given pair  $\langle \text{English parse tree } e, \text{ foreign language string } f \rangle$ , this model returns a probability  $P(f|e)$ .
- A Language model (LM). For any English tree  $e$ , this model returns a probability  $P(e)$ .
- A Search algorithm. Given a foreign language sentence  $f$ , this algorithm searches for the English tree  $e$  that maximizes  $P(e|f) \approx P(e) \cdot P(f|e)$ .

Yamada and Knight (2001; 2002) describe the TM and the search algorithm, respectively. The search algorithm takes a foreign language sentence and produces a vast number of candidate English trees, packed into a forest structure, as in summarization (see Section 2), then searches for the highest-scoring tree. The current algorithm uses a trigram LM, ignoring the internal structure of the

candidate trees. Therefore, the system does not necessarily produce syntactically correct translations.

To see the effect of a lexicalized syntax language model, we performed what automatic speech recognition (ASR) researchers call “a cheating experiment”. For a given acoustic signal, ASR researchers know both the correct target transcription, A, (which was done by hand) and the current automatic system transcription, B. The probabilistic score for B will be greater than that for A. However, a new knowledge source may provide additional scores that cause A to be reranked higher. This is cheating for two reasons: (1) it does not include a search algorithm that integrates the new knowledge source, and (2) there may be an incorrect string C that scores higher than both A and B under reranking.

This kind of experiment is not regularly done in machine translation because there is no single correct translation A. Using a human translation as a target A does not work, because human translations are often non-literal and current statistical models do not recognize them as good translations. To remedy this, we manually created a set of target translations, which we called “hope” translations. These are good translations that we believe to be within reach of the system: we can reasonably hope that the system would prefer them.

In our experiment, we score both hope sentences, A, and current system translations, B, over a number of examples, using combinations of these knowledge sources:

- T: translation model
- R: word-trigram language model
- C: score from Charniak’s parser

Table 5 shows the results. A sentence marked “dec1” is a decoder output (Sentence B) and “hope” is a hope sentence (Sentence A); lower scores are better. The last row shows the average difference of the score (a positive difference means that the system prefers “hope” sentences over current system outputs). Just above the last row is the number of sentences which ranked better.<sup>3</sup>

<sup>3</sup>The score from Charniak’s parser (C) is a  $-\log(\text{un-normalized prob})$ , thus it may yield negative value. The TM scores are calculated from parse trees, not from

			T	R	C	T+R	T+2R	T+C	T+R+C	T+2C	T+R+2C
1	decl hope	great attention his first visit to china his first visit is highly important to the chinese side	36.36 66.74	39.63 54.17	17.44 -5.14	75.99 120.91	115.62 175.08	53.80 61.60	93.43 115.77	71.24 56.46	110.87 110.63
2	decl hope	he briefed reporters declaring the major contents he declared the main contents to the press agency	45.47 64.81	52.62 56.69	16.53 0.72	98.09 121.50	150.72 178.20	62.00 65.52	114.62 122.22	78.53 66.24	131.15 122.93
3	decl hope	industrial production and marketing join basic level restored to normal state the basic link between supply and marketing in industry resumed normal status	50.84 81.84	74.34 83.53	25.62 15.01	125.17 165.37	199.51 248.90	76.45 96.85	150.79 180.38	102.07 111.86	176.41 195.39
4	decl hope	reminbi exchange rate stability this reputation the reminbi exchange rate continued to maintain stability	16.99 27.34	43.38 45.43	36.31 18.27	60.37 72.77	103.75 118.20	53.30 45.62	96.68 91.05	89.61 63.89	132.99 109.32
5	decl hope	order form is really "agricultural" the rapid development third. " agriculture to order " has developed rapidly	27.32 47.61	57.36 62.06	25.96 17.79	84.68 109.67	142.04 171.73	53.29 65.40	110.65 127.46	79.25 83.19	136.61 145.25
6	decl hope	talks have topic hot spot talks will have many hot topics	23.65 39.73	49.54 47.08	28.26 13.31	73.19 86.81	122.73 133.88	51.91 53.03	101.45 100.11	80.17 66.34	129.71 113.42
7	decl hope	reminding one person alone can not treat to remind chen that one-man rule is unacceptable	39.80 44.26	42.88 73.31	16.20 39.15	82.69 117.56	125.57 190.87	56.01 83.41	98.89 156.71	72.21 122.55	115.09 195.86
8	decl hope	bad now like this ? now it is not good ?	45.15 17.69	30.03 32.87	7.01 34.18	75.18 50.56	105.21 83.43	52.16 51.87	82.19 84.74	59.17 86.05	89.21 118.92
9	decl hope	stratify binomthly the masses the masses are stratified	26.52 17.69	34.07 32.87	11.30 34.18	60.58 50.56	94.65 83.43	57.82 51.87	71.89 84.74	49.12 86.05	83.19 118.92
10	decl hope	undercompensation conditions for accept the us side to this aspect the us side did not accept the dprk side 's proposed compensation terms	51.95 52.35	53.75 70.89	38.74 27.98	105.70 123.24	159.45 194.13	90.69 80.33	144.44 151.22	129.43 108.31	183.18 179.19
11	decl hope	this is extremely absurd this is extremely absurd	21.93 24.15	18.81 18.81	2.01 2.01	40.74 42.96	59.54 61.76	23.94 26.16	42.74 44.96	25.95 28.16	44.75 46.97
12	decl hope	rights and obligations are associated with each other but rights and obligations are mutually linked	40.96 35.73	32.37 52.06	-3.53 12.32	73.33 87.79	105.69 139.84	37.44 48.05	69.80 100.10	33.91 60.37	66.27 112.42
13	decl hope	mass various channels to reflect the views there are many such channels to reflect the views of the masses	39.88 57.78	39.22 45.56	16.68 -4.22	79.10 103.33	118.32 148.89	56.55 53.56	95.77 99.11	73.23 49.34	112.45 94.89
14	decl hope	hong kong is our home hong kong is our home	13.93 16.78	23.39 23.39	12.34 12.34	37.32 40.17	60.71 63.55	26.27 29.12	49.66 52.30	38.61 41.45	61.99 64.84
15	decl hope	he said china evermore taste he said he often eats chinese food	31.68 47.46	32.81 54.21	27.96 8.63	64.49 101.67	97.31 155.88	59.64 56.09	92.45 110.30	87.59 64.72	120.41 118.93
16	decl hope	hu jintao said hu jintao said	5.52 7.90	14.57 14.57	39.17 39.17	20.09 22.47	34.66 37.04	44.69 47.07	59.26 61.64	83.86 86.24	98.43 100.81
17	decl hope	this is not possible --this is not possible	22.03 21.26	17.21 19.97	-3.10 4.00	39.24 41.23	56.44 61.20	18.93 25.26	36.14 45.23	15.83 29.27	33.04 49.24
18	decl hope	this globalization of production this is the globalization of production	18.67 22.08	25.89 25.90	10.68 -0.59	44.56 47.98	70.46 73.88	29.35 21.49	55.24 47.39	40.03 20.90	65.92 46.80
19	decl hope	many successes possible ? how much success is possible ?	24.15 28.77	32.48 35.91	13.53 -4.54	56.63 56.63	89.11 100.59	37.68 24.23	70.16 60.14	51.21 19.70	83.69 55.61
20	decl hope	clinton meeting was held this idea holding this meeting was clinton 's decision	42.07 52.66	40.63 49.59	15.04 14.71	82.69 102.24	123.32 151.83	57.11 67.37	97.73 116.95	72.15 82.08	112.78 131.66
21	decl hope	reiterative he explained he has repeatedly explained	10.10 18.72	17.77 21.66	23.17 7.40	27.87 40.39	45.63 62.05	33.27 26.12	51.03 47.78	36.44 33.52	74.20 55.18
T	decl hope		18 3	19 5	6 18	20 1	20 1	12 9	14 7	9 12	9 12
ave. diff			-1.048	-0.776	0.956	-1.825	-2.601	-0.092	-0.868	0.864	0.088

Table 5: Results of the experiment on the MT system

```
(H37 :ADJUNCT "earlier"
      :LOGICAL-SUBJECT (H5 / "company")
      / "announce"
      :ADJUNCT (H34 :ADJUNCT "its"
                    / "plan"))
```

Figure 1: An underspecified input for the sentence, "Earlier the company announced its plans."

As expected, our current system (T+R) almost never ranks the hope sentence higher than the system sentence. If we replace the trigram model with Charniak’s parser (T+C), the hope sentence is ranked higher than the originally preferred sentences in 9 cases. These results are slightly better than reranking with the parser (T+R+C). Best results are obtained by assigning a higher weight to the parser score relative to the translation model (T+2C): here, the hope sentence comes out on top in 12 cases.

## 4 Natural Language Generation

From the sentences in section 23 of the Penn Treebank, inputs to the HALogen generator system were automatically derived and then regenerated (Langkilde-Geary, 2002). The inputs derived were feature-value dependency structures, where the features represent syntactic relationships between values, and the values were either words in root form or a nested feature-value structure. The inputs were underspecified with respect to properties such as part-of-speech category, tense, voice, and number, as well as constituent order and some closed-class words like auxiliary verbs and determiners. Underspecification tests the ability of a language model to pick the best solution from among a set of choices.

HALogen overgenerates possible expressions for an input, in part because of enumerating possible choices for underspecified details. It then ranks potential outputs using an  $n$ -gram model. Both bigram and trigram models were available, but because the trigram model takes two orders of magnitude more time to use (roughly 40 minutes per sentence, versus 1 minute per sentence for the bigram model), the sentences were generated using the bigram model.

sentences. For “decl” sentences, we use the parse tree returned from the decoder to calculate it. For “hope” sentences, we use a parse tree generated from Collins parser (1997) to calculate it. Thus, there may be different T scores for the same sentence.

One hundred sentences were randomly chosen from among the successfully generated outputs. Each was paired with its original Treebank sentence, and then trigram and Charniak-parser scores were calculated for all the sentences. About 3% of the sentence pairs were exact matches. Sentences within pairs had very similar lengths in all cases. The average sentence length was 24 tokens. The original Treebank sentence serves as a gold standard for the generator.

Trigram scores for original Treebank sentences scored better than the output of the generator system 71% of the time. In comparison, Charniak-parse scores preferred the original Treebank sentences 83% of the time. This indicates that the generator system would benefit even more from a statistical model of syntax than from trigrams.

## 5 Conclusion

In statistical summarization, we have shown that having a lexicalized syntax model reduces the frequency of grammatical errors through both a careful error analysis and a human evaluation. We furthermore show that a lexicalized syntax model might assist in the selection of good translations in a syntax-based machine translation system. Finally, we present results that indicate the importance of lexicalized models of syntax in the HALogen natural language generation system (Langkilde-Geary, 2002).

It seems from these results that the integration of such a language model in these tasks and in statistical natural language generation systems will prove to be fruitful. Most importantly, all three of these systems use the same forest ranking algorithm/component. Thus, if this one component were extended to use a lexicalized model of syntax in place of the current  $n$ -gram scoring method, the performance of all three system would undoubtedly improve significantly.

## Acknowledgements

This work was partially supported by DARPA-ITO grant N66001-00-1-9814 and a USC Dean Fellowship to Hal Daume III.



## References

- S. Bangalore and O. Rambow. 2000. Exploiting a probabilistic hierarchical model for generation. In *Proceedings of the International Conference on Computational Linguistics (COLING 2000)*, Hong Kong, China.
- Eugene Charniak. 1997. Statistical parsing with a context-free grammar and word statistics. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence (AAAI'97)*, pages 598–603, Providence, Rhode Island, July 27–31.
- Eugene Charniak. 2000. A maximum-entropy-inspired parser. In *Proceedings of the First Annual Meeting of the North American Chapter of the Association for Computational Linguistics NAACL-2000*, pages 132–139, Seattle, Washington, April 29 – May 3.
- Eugene Charniak. 2001. Immediate-head parsing for language models. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*.
- Ciprian Chelba and Frederick Jelinek. 1998. Exploiting syntactic structure for language modeling. In Christian Boitet and Pete Whitelock, editors, *Proceedings of the Thirty-Sixth Annual Meeting of the Association for Computational Linguistics and Seventeenth International Conference on Computational Linguistics*, pages 225–231, San Francisco, California. Morgan Kaufmann Publishers.
- Michael Collins. 1997. Three generative, lexicalized models for statistical parsing. In *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (ACL-97)*, pages 16–23, Madrid, Spain, July 7–12.
- Hal Daumé III and Daniel Marcu. 2002. A noisy-channel model for document compression. In *Proceedings of the Conference of the Association of Computational Linguistics (ACL 2002)*.
- L. Hirschman, M. Light, E. Breck, and J. Burger. 1999. Deep read: A reading comprehension system. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*.
- Kevin Knight and Daniel Marcu. 2000. Statistics-based summarization — step one: Sentence compression. In *The 17th National Conference on Artificial Intelligence (AAAI-2000)*, pages 703–710, Austin, TX, July 30th – August 3rd.
- Irene Langkilde and Kevin Knight. 1998. Generation that exploits corpus-based statistical knowledge. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and of the 17th International Conference on Computational Linguistics (COLING/ACL'98)*, Montreal, Canada, August.
- Irene Langkilde-Geary. 2002. An empirical verification of coverage and correctness for a general-purpose sentence generator. In *Proceedings of the 2nd International Conference on Natural Language Generation (INLG 2002)*., Arden House, NJ, July.
- Irene Langkilde. 2000. Forest-based statistical sentence generation. In *Proceedings of the 1st Annual Meeting of the North American Chapter of the Association for Computational Linguistics*, Seattle, Washington, April 30–May 3.
- Brian Roark. 2001. Probabilistic top-down parsing and language modelling. In *Computational Linguistics* 27(2), pages 249–276.
- K. Yamada and K. Knight. 2001. A syntax-based statistical translation model. In *Proceedings of the Conference of the Association of Computational Linguistics (ACL 2001)*.
- K. Yamada and K. Knight. 2002. A decoder for syntax-based statistical MT. In *Proceedings of the Conference of the Association of Computational Linguistics (ACL 2002)*.